**Activity**: Hands-on Activity 7.2 Webscraping using BeautifulSoup and Requests

**Name**: Juliann Vincent B. Quibral

**Section**: BSCPE22S3

## ⌄ Example of gathering image data using webcam

**Open Data and Private Data**

1. Open Data The Open Knowledge Foundation describes Open Data as "any content, information or data that people are free to use, reuse, and redistribute without any legal, technological, or social restriction."
2. Private Data Data related to an expectation of privacy and regulated by a particular country/government

**Structured and Unstructured Data**

1. Structured Data Data entered and maintained in fixed fields within a file or record Easily entered, classified, queried, and analyzed Relational databases or spreadsheets
2. Unstructured Data Lacks organization Raw data Photo contents, audio, video, web pages, blogs, books, journals, white papers, PowerPoint presentations, articles, email, wikis, word processing documents, and text in general

```python
1  import cv2
2  from google.colab.patches import cv2_imshow
3
4  webcam = cv2.VideoCapture(0)
5
6  while True:
7      try:
8          check, frame = webcam.read()
9          if check:  # Check if frame is successfully captured
10             print(check)  # prints true as long as the webcam is running
11             print(frame)  # prints matrix values of each frame
12             cv2_imshow(frame)  # Display the frame
13             key = cv2.waitKey(1)
14
15             if key == ord('s'):
16                 cv2.imwrite(filename='saved_img.jpg', img=frame)
17                 webcam.release()
18                 img_new = cv2.imread('saved_img.jpg', cv2.IMREAD_GRAYSCALE)
19                 img_new = cv2_imshow(img_new)
20                 cv2.waitKey(1650)
21                 cv2.destroyAllWindows()
22                 print("Processing image...")
23                 img_ = cv2.imread('saved_img.jpg', cv2.IMREAD_ANYCOLOR)
24                 print("Converting RGB image to grayscale...")
25                 gray = cv2.cvtColor(img_, cv2.COLOR_BGR2GRAY)
26                 print("Converted RGB image to grayscale...")
27                 print("Resizing image to 28x28 scale...")
28                 img_ = cv2.resize(gray, (28, 28))
29                 print("Resized...")
30                 img_resized = cv2.imwrite(filename='saved_img-final.jpg', img=img_)
31                 print("Image saved!")
32                 break
33
34             elif key == ord('q'):
35                 print("Turning off camera.")
36                 webcam.release()
37                 print("Camera off.")
38                 print("Program ended.")
39                 cv2.destroyAllWindows()
40                 break
41         else:
42             print("Unable to capture frame. Check your webcam connection.")
43             break
44
45     except KeyboardInterrupt:
46         print("Turning off camera.")
47         webcam.release()
48         print("Camera off.")
49         print("Program ended.")
```

```
50        cv2.destroyAllWindows()
51        break
52
```

```
Unable to capture frame. Check your webcam connection.
```

## ∨  Example of gathering voice data using microphone

```
1 !pip3 install sounddevice
```

```
Collecting sounddevice
  Downloading sounddevice-0.4.6-py3-none-any.whl (31 kB)
Requirement already satisfied: CFFI>=1.0 in /usr/local/lib/python3.10/dist-packages (from sounddevice) (1.16.0)
Requirement already satisfied: pycparser in /usr/local/lib/python3.10/dist-packages (from CFFI>=1.0->sounddevice) (2.21)
Installing collected packages: sounddevice
Successfully installed sounddevice-0.4.6
```

```
1 !pip3 install wavio
```

```
Collecting wavio
  Downloading wavio-0.0.8-py3-none-any.whl (9.4 kB)
Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.10/dist-packages (from wavio) (1.25.2)
Installing collected packages: wavio
Successfully installed wavio-0.0.8
```

```
1 !pip3 install scipy
```

```
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (1.11.4)
Requirement already satisfied: numpy<1.28.0,>=1.21.6 in /usr/local/lib/python3.10/dist-packages (from scipy) (1.25.2)
```

```
1 !apt-get install libportaudio2
```

```
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following NEW packages will be installed:
  libportaudio2
0 upgraded, 1 newly installed, 0 to remove and 39 not upgraded.
Need to get 65.3 kB of archives.
After this operation, 223 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libportaudio2 amd64 19.6.0-1.1 [65.3 kB]
Fetched 65.3 kB in 1s (81.0 kB/s)
Selecting previously unselected package libportaudio2:amd64.
(Reading database ... 121753 files and directories currently installed.)
Preparing to unpack .../libportaudio2_19.6.0-1.1_amd64.deb ...
Unpacking libportaudio2:amd64 (19.6.0-1.1) ...
Setting up libportaudio2:amd64 (19.6.0-1.1) ...
Processing triggers for libc-bin (2.35-0ubuntu3.4) ...
/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_5.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbb.so.12 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc_proxy.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_0.so.3 is not a symbolic link
```

```
1 !pip install sounddevice --upgrade
```

```
Requirement already satisfied: sounddevice in /usr/local/lib/python3.10/dist-packages (0.4.6)
Requirement already satisfied: CFFI>=1.0 in /usr/local/lib/python3.10/dist-packages (from sounddevice) (1.16.0)
Requirement already satisfied: pycparser in /usr/local/lib/python3.10/dist-packages (from CFFI>=1.0->sounddevice) (2.21)
```

```
 1 # import required libraries
 2 import sounddevice as sd
 3 from scipy.io.wavfile import write
 4 import wavio as wv
 5
 6 # Sampling frequency
 7 freq = 48000
 8
 9 # Recording duration
10 duration = 5
11
12 # Start recorder with the given values
13 # of duration and sample frequency
14 recording = sd.rec(int(duration * freq),
15  samplerate=freq, channels=2)
16
17 # Record audio for the given number of seconds
18 sd.wait()
19
20 # This will convert the NumPy array to an audio
21 # file with the given sampling frequency
22 write("recording0.wav", freq, recording)
23 # Convert the NumPy array to audio file
24 wv.write("recording1.wav", recording, freq, sampwidth=2)
```

```
---------------------------------------------------------------------------
PortAudioError                            Traceback (most recent call last)
<ipython-input-21-3f46ebebbb4e> in <cell line: 14>()
     12 # Start recorder with the given values
     13 # of duration and sample frequency
---> 14 recording = sd.rec(int(duration * freq),
     15  samplerate=freq, channels=2)
     16

                          ⇕ 5 frames
/usr/local/lib/python3.10/dist-packages/sounddevice.py in query_devices(device, kind)
    567     info = _lib.Pa_GetDeviceInfo(device)
    568     if not info:
--> 569         raise PortAudioError(f'Error querying device {device}')
    570     assert info.structVersion == 2
    571     name_bytes = _ffi.string(info.name)

PortAudioError: Error querying device -1
```

## Web Scraping

```
1 !pip install bs4
```

```
Collecting bs4
  Downloading bs4-0.0.2-py2.py3-none-any.whl (1.2 kB)
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.10/dist-packages (from bs4) (4.12.3)
Requirement already satisfied: soupsieve>1.2 in /usr/local/lib/python3.10/dist-packages (from beautifulsoup4->bs4) (2.5)
Installing collected packages: bs4
Successfully installed bs4-0.0.2
```

```
1 pip install requests
```

```
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests) (3.6)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests) (2024.2.2)
```

```python
 1 import requests
 2 from bs4 import BeautifulSoup
 3
 4 def getdata(url):
 5   r = requests.get(url)
 6   return r.text
 7
 8 htmldata = getdata("https://www.google.com/")
 9 soup = BeautifulSoup(htmldata, 'html.parser')
10 for item in soup.find_all('img'):
11   print(item['src'])
12
```

```
/images/branding/googlelogo/1x/googlelogo_white_background_color_272x92dp.png
```

```
1 pip install selenium
```

```
Collecting selenium
  Downloading selenium-4.18.1-py3-none-any.whl (10.0 MB)
     ──────────────────────────────────────── 10.0/10.0 MB 52.9 MB/s eta 0:00:00
Requirement already satisfied: urllib3[socks]<3,>=1.26 in /usr/local/lib/python3.10/dist-packages (from selenium) (2.0.7)
Collecting trio~=0.17 (from selenium)
  Downloading trio-0.25.0-py3-none-any.whl (467 kB)
     ──────────────────────────────────────── 467.2/467.2 kB 40.4 MB/s eta 0:00:00
Collecting trio-websocket~=0.9 (from selenium)
  Downloading trio_websocket-0.11.1-py3-none-any.whl (17 kB)
Requirement already satisfied: certifi>=2021.10.8 in /usr/local/lib/python3.10/dist-packages (from selenium) (2024.2.2)
Requirement already satisfied: typing_extensions>=4.9.0 in /usr/local/lib/python3.10/dist-packages (from selenium) (4.10.0)
Requirement already satisfied: attrs>=23.2.0 in /usr/local/lib/python3.10/dist-packages (from trio~=0.17->selenium) (23.2.0)
Requirement already satisfied: sortedcontainers in /usr/local/lib/python3.10/dist-packages (from trio~=0.17->selenium) (2.4.0)
Requirement already satisfied: idna in /usr/local/lib/python3.10/dist-packages (from trio~=0.17->selenium) (3.6)
Collecting outcome (from trio~=0.17->selenium)
  Downloading outcome-1.3.0.post0-py2.py3-none-any.whl (10 kB)
Requirement already satisfied: sniffio>=1.3.0 in /usr/local/lib/python3.10/dist-packages (from trio~=0.17->selenium) (1.3.1)
Requirement already satisfied: exceptiongroup in /usr/local/lib/python3.10/dist-packages (from trio~=0.17->selenium) (1.2.0)
Collecting wsproto>=0.14 (from trio-websocket~=0.9->selenium)
  Downloading wsproto-1.2.0-py3-none-any.whl (24 kB)
Requirement already satisfied: pysocks!=1.5.7,<2.0,>=1.5.6 in /usr/local/lib/python3.10/dist-packages (from urllib3[socks]<3,>=1.26->sel
Collecting h11<1,>=0.9.0 (from wsproto>=0.14->trio-websocket~=0.9->selenium)
  Downloading h11-0.14.0-py3-none-any.whl (58 kB)
     ──────────────────────────────────────── 58.3/58.3 kB 8.2 MB/s eta 0:00:00
Installing collected packages: outcome, h11, wsproto, trio, trio-websocket, selenium
Successfully installed h11-0.14.0 outcome-1.3.0.post0 selenium-4.18.1 trio-0.25.0 trio-websocket-0.11.1 wsproto-1.2.0
```

## Image Scraping using Selenium

```python
1  !pip install selenium
2  !apt-get update # to update ubuntu to correctly run apt install
3  !apt install chromium-chromedriver
4  !cp /usr/lib/chromium-browser/chromedriver /usr/bin
5  import sys
6  sys.path.insert(0,'/usr/lib/chromium-browser/chromedriver')
7
8  from selenium import webdriver
9  import time
10 import requests
11 import shutil
12 import os
13 import getpass
14 import urllib.request
15 import io
16 import time
17 from PIL import Image
18 user = getpass.getuser()
19 chrome_options = webdriver.ChromeOptions()
20 chrome_options.add_argument('--headless')
21 chrome_options.add_argument('--no-sandbox')
22 chrome_options.add_argument('--disable-dev-shm-usage')
23 driver = webdriver.Chrome('chromedriver',chrome_options=chrome_options)
24
25 search_url = "https://www.google.com/search?q={q}&tbm=isch&tbs=sur%3Afc&hl=en&ved=0CAIQpwVqFwoTCKCa1c6s4-oCFQAAAAAdAAAAABAC&biw=1251&bih=
26 driver.get(search_url.format(q='Car'))
27
28 def scroll_to_end(driver):
29   driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
30   time.sleep(5)#sleep_between_interactions
31
32 def getImageUrls(name,totalImgs,driver):
33   search_url = "https://www.google.com/search?q={q}&tbm=isch&tbs=sur%3Afc&hl=en&ved=0CAIQpwVqFwoTCKCa1c6s4-oCFQAAAAAdAAAAABAC&biw=1251&bi
34   driver.get(search_url.format(q=name))
35   img_urls = set()
36   img_count = 0
37   results_start = 0
38
39   while(img_count<totalImgs): #Extract actual images now
40
41     scroll_to_end(driver)
42
43     thumbnail_results = driver.find_elements_by_xpath("//img[contains(@class,'Q4LuWd')]")
44     totalResults=len(thumbnail_results)
45     print(f"Found: {totalResults} search results. Extracting links from{results_start}:{totalResults}")
46
47     for img in thumbnail_results[results_start:totalResults]:
48
49       img.click()
50       time.sleep(2)
51       actual_images = driver.find_elements_by_css_selector('img.n3VNCb')
52       for actual_image in actual_images:
53         if actual_image.get_attribute('src') and 'https' in actual_image.get_attribute('src'):
54           img_urls.add(actual_image.get_attribute('src'))
55
56       img_count=len(img_urls)
57
58       if img_count >= totalImgs:
59         print(f"Found: {img_count} image links")
60         break
61       else:
62         print("Found:", img_count, "looking for more image links ...")
63         load_more_button = driver.find_element_by_css_selector(".mye4qd")
64         driver.execute_script("document.querySelector('.mye4qd').click();")
65         results_start = len(thumbnail_results)
66   return img_urls
67
68 def downloadImages(folder_path,file_name,url):
69     try:
70       image_content = requests.get(url).content
71     except Exception as e:
72       print(f"ERROR - COULD NOT DOWNLOAD {url} - {e}")
73
74     try:
75       image_file = io.BytesIO(image_content)
76       image = Image.open(image_file).convert('RGB')
77
```

```python
78           file_path = os.path.join(folder_path, file_name)
79
80           with open(file_path, 'wb') as f:
81             image.save(f, "JPEG", quality=85)
82           print(f"SAVED - {url} - AT: {file_path}")
83       except Exception as e:
84         print(f"ERROR - COULD NOT SAVE {url} - {e}")
85
86   def saveInDestFolder(searchNames,destDir,totalImgs,driver):
87       for name in list(searchNames):
88         path=os.path.join(destDir,name)
89         if not os.path.isdir(path):
90           os.mkdir(path)
91         print('Current Path',path)
92         totalLinks=getImageUrls(name,totalImgs,driver)
93         print('totalLinks',totalLinks)
94       if totalLinks is None:
95           print('images not found for :',name)
96
97       else:
98         for i, link in enumerate(totalLinks):
99           file_name = f"{i:150}.jpg"
100          downloadImages(path,file_name,link)
101
102  searchNames=['cat']
103  destDir=f'/content/drive/My Drive/Colab Notebooks/Dataset/'
104  totalImgs=5
105
106  saveInDestFolder(searchNames,destDir,totalImgs,driver)
```

```
Requirement already satisfied: selenium in /usr/local/lib/python3.10/dist-packages (4.18
Requirement already satisfied: urllib3[socks]<3,>=1.26 in /usr/local/lib/python3.10/dist
Requirement already satisfied: trio~=0.17 in /usr/local/lib/python3.10/dist-packages (fr
Requirement already satisfied: trio-websocket~=0.9 in /usr/local/lib/python3.10/dist-pac
Requirement already satisfied: certifi>=2021.10.8 in /usr/local/lib/python3.10/dist-pack
Requirement already satisfied: typing_extensions>=4.9.0 in /usr/local/lib/python3.10/dis
Requirement already satisfied: attrs>=23.2.0 in /usr/local/lib/python3.10/dist-packages
Requirement already satisfied: sortedcontainers in /usr/local/lib/python3.10/dist-packag
Requirement already satisfied: idna in /usr/local/lib/python3.10/dist-packages (from tri
Requirement already satisfied: outcome in /usr/local/lib/python3.10/dist-packages (from
Requirement already satisfied: sniffio>=1.3.0 in /usr/local/lib/python3.10/dist-packages
Requirement already satisfied: exceptiongroup in /usr/local/lib/python3.10/dist-packages
Requirement already satisfied: wsproto>=0.14 in /usr/local/lib/python3.10/dist-packages
Requirement already satisfied: pysocks!=1.5.7,<2.0,>=1.5.6 in /usr/local/lib/python3.10/
Requirement already satisfied: h11<1,>=0.9.0 in /usr/local/lib/python3.10/dist-packages
Get:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease [3,626 B]
Get:2 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64  InRele
Get:3 http://security.ubuntu.com/ubuntu jammy-security InRelease [110 kB]
Hit:4 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:5 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64  Packag
Get:6 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [119 kB]
Hit:7 http://archive.ubuntu.com/ubuntu jammy-backports InRelease
Hit:8 https://ppa.launchpadcontent.net/c2d4u.team/c2d4u4.0+/ubuntu jammy InRelease
Get:9 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages [1,356 kB]
Hit:10 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Hit:11 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy InRelease
Get:12 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [1,898 kB]
Hit:13 https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu jammy InRelease
Fetched 4,261 kB in 2s (2,625 kB/s)
Reading package lists... Done
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  apparmor chromium-browser libfuse3-3 liblzo2-2 libudev1 snapd squashfs-tools systemd-h
  udev
Suggested packages:
  apparmor-profiles-extra apparmor-utils fuse3 zenity | kdialog
The following NEW packages will be installed:
  apparmor chromium-browser chromium-chromedriver libfuse3-3 liblzo2-2 snapd squashfs-to
  systemd-hwe-hwdb udev
The following packages will be upgraded:
  libudev1
1 upgraded, 9 newly installed, 0 to remove and 38 not upgraded.
Need to get 26.4 MB of archives.
After this operation, 116 MB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 apparmor amd64 3.0.4-2ub
Get:2 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 liblzo2-2 amd64 2.10-2build3 [5:
Get:3 http://archive.ubuntu.com/ubuntu jammy/main amd64 squashfs-tools amd64 1:4.5-3buil
Get:4 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libudev1 amd64 249.11-0u
Get:5 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 udev amd64 249.11-0ubunt
Get:6 http://archive.ubuntu.com/ubuntu jammy/main amd64 libfuse3-3 amd64 3.10.5-1build1
Get:7 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 snapd amd64 2.58+22.04.1
Get:8 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 chromium-browser amd
Get:9 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 chromium-chromedrive
Get:10 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 systemd-hwe-hwdb all 24
Fetched 26.4 MB in 1s (18.6 MB/s)
```

## Web Scraping of Movies Information using BeautifulSoup

```python
1 from requests import get
2 url = 'https://hurawatch.cc/search/logan'
3 response = get(url)
4 print(response.text[:500])
```

```
<!DOCTYPE html>
<html xmlns="http://www.w3.org/1999/xhtml" xml:lang="en" lang="en">
<head>
    <title>Search results for &#39;logan&#39; Movies &amp; Tv Series Hurawatch</title>
<meta http-equiv="Content-Type" content="text/html; charset=UTF-8"/>

    <meta name="robots" content="index, follow">
    <meta name="revisit-after" content="1 days">

<meta http-equiv="content-language" content="en"/>
<link rel="dns-prefetch" href="//www.google-analytics.com">
<link rel="dns-prefetch" href="//www.gstat
```

```
1 from bs4 import BeautifulSoup
2 html_soup = BeautifulSoup(response.text, 'html.parser')
3 headers = {'Accept-Language': 'en-US,en;q=0.8'}
4 type(html_soup)
5
```

```
bs4.BeautifulSoup
def __call__(*args, **kwargs)

A data structure representing a parsed HTML or XML document.

Most of the methods you'll call on a BeautifulSoup object are inherited from
PageElement or Tag.

Internally, this class defines the basic interface called by the
```

```
1 movie_containers = html_soup.find_all('div', class_ = 'flw-item')
2 print(type(movie_containers))
3 print(len(movie_containers))
```

```
    <class 'bs4.element.ResultSet'>
    6
```

```
1 first_movie = movie_containers[0]
2 first_movie
```

```
    <div class="flw-item">
    <div class="film-poster">
    <div class="pick film-poster-quality">HD</div>
    <img alt="Logan" class="film-poster-img lazyload" data-
    src="https://img.hurawatch.cc/xxrz/250x400/348/16/f7/16f7b48a3df281f25cb746394488ea9d/16f7b48a3df281f25cb746394488ea9d.jpg"
    title="Logan"/>
    <a class="film-poster-ahref flw-item-tip" href="/movie/watch-logan-online-19754" title="Logan"><i class="fa fa-play"></i></a>
    </div>
    <div class="film-detail">
    <h2 class="film-name"><a href="/movie/watch-logan-online-19754" title="Logan">Logan</a>
    </h2>
    <div class="fd-infor">
    <span class="fdi-item">2017</span>
    <span class="dot"></span>
    <span class="fdi-item fdi-duration">137m</span>
    <span class="float-right fdi-type">Movie</span>
    </div>
    <div class="clearfix"></div>
    </div>
    <div class="clearfix"></div>
    </div>
```

```
1 first_movie.div
```

```
    <div class="film-poster">
    <div class="pick film-poster-quality">HD</div>
    <img alt="Logan" class="film-poster-img lazyload" data-
    src="https://img.hurawatch.cc/xxrz/250x400/348/16/f7/16f7b48a3df281f25cb746394488ea9d/16f7b48a3df281f25cb746394488ea9d.jpg"
    title="Logan"/>
    <a class="film-poster-ahref flw-item-tip" href="/movie/watch-logan-online-19754" title="Logan"><i class="fa fa-play"></i></a>
    </div>
```

```
1 first_movie.a
```

```
    <a class="film-poster-ahref flw-item-tip" href="/movie/watch-logan-online-19754" title="Logan"><i class="fa fa-play"></i></a>
```

```
1 first_movie.h2
```

```
    <h2 class="film-name"><a href="/movie/watch-logan-online-19754" title="Logan">Logan</a>
    </h2>
```

```
1 first_movie.h2.a
```

```
    <a href="/movie/watch-logan-online-19754" title="Logan">Logan</a>
```

```
1 first_name = first_movie.h2.a.text
2 first_name
```

```
    'Logan'
```

```
1 first_year = first_movie.find('span', class_='fdi-item')
2 if first_year:
3     print(first_year.text)
4 else:
5     print("Year information not found")
6
```

```
    2017
```

```
1 first_year = first_year.text
2 first_year
```

```
    '2017'
```

Rating doesn't exist in my link

```
1 first_movie.strong
```

```
1 first_imdb = float(first_movie.strong.text)
2 first_imdb
```

```
    ---------------------------------------------------------------------------
    AttributeError                            Traceback (most recent call last)
    <ipython-input-88-92faeb51c9f2> in <cell line: 1>()
    ----> 1 first_imdb = float(first_movie.strong.text)
          2 first_imdb

    AttributeError: 'NoneType' object has no attribute 'text'
```

Metascore doesn't exist in my link

```
1 first_mscore = first_movie.find('span', class_ = 'metascore favorable')
2 first_mscore = int(first_mscore.text)
3 print(first_mscore)
4
```

```
    ---------------------------------------------------------------------------
    AttributeError                            Traceback (most recent call last)
    <ipython-input-89-889bc009bd72> in <cell line: 2>()
          1 first_mscore = first_movie.find('span', class_ = 'metascore favorable')
    ----> 2 first_mscore = int(first_mscore.text)
          3 print(first_mscore)

    AttributeError: 'NoneType' object has no attribute 'text'
```

Votes doesn't exist in my link

```
1 first_votes = first_movie.find('span', attrs = {'name':'nv'})
2 first_votes
```

```
1 first_votes['data-value']
```

```
    ---------------------------------------------------------------------------
    TypeError                                 Traceback (most recent call last)
    <ipython-input-92-2d836d02a09a> in <cell line: 1>()
    ----> 1 first_votes['data-value']

    TypeError: 'NoneType' object is not subscriptable
```

```
1 first_votes = int(first_votes['data-value'])
```

```
-----------------------------------------------------------------------
TypeError                                 Traceback (most recent call last)
<ipython-input-93-e337b21fe258> in <cell line: 1>()
----> 1 first_votes = int(first_votes['data-value'])

TypeError: 'NoneType' object is not subscriptable
```

However the site have a time duration

```
1 first_duration = first_movie.find('span', class_='fdi-item fdi-duration')
2 if first_duration:
3     print(first_duration.text)
4 else:
5     print("Duration information not found")
6
```

```
137m
```

```
1 # Lists to store the scraped data in
2 names = []
3 years = []
4 durations = []
5
6
7 for container in movie_containers:
8     if container.find('div', class_='fd-infor') is not None:
9         # Name
10        name = container.h2.a.text
11        names.append(name)
12
13        # Year
14        year = container.find('span', class_='fdi-item').text
15        years.append(year)
16
17        # Duration
18        duration_element = container.find('span', class_='fdi-item fdi-duration')
19        if duration_element is not None:
20            duration = duration_element.text
21        else:
22            duration = 'Not available'
23        durations.append(duration)
24
25 print(names)
26 print(years)
27 print(durations)
```

```
['Logan', 'Logan Lucky', 'The Night Logan Woke Up', 'The Taking of Deborah Logan', 'The Two Worlds of Jennie Logan', "Logan's Run"]
['2017', '2017', 'SS 1', '2014', '1979', '1976']
['137m', '119m', 'Not available', '90m', '94m', '119m']
```

```
1 import pandas as pd
2 test_df = pd.DataFrame({
3 'movie': names,
4 'year': years,
5 'duration': duration,
6 })
7 print(test_df.info())
8 test_df
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   movie     6 non-null      object
 1   year      6 non-null      object
 2   duration  6 non-null      object
dtypes: object(3)
memory usage: 272.0+ bytes
None
```

|   | movie | year | duration |
|---|-------|------|----------|
| 0 | Logan | 2017 | 119m |
| 1 | Logan Lucky | 2017 | 119m |
| 2 | The Night Logan Woke Up | SS 1 | 119m |
| 3 | The Taking of Deborah Logan | 2014 | 119m |
| 4 | The Two Worlds of Jennie Logan | 1979 | 119m |
| 5 | Logan's Run | 1976 | 119m |

My website doesn't have another batch of movies and is limited to just 6/6

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype
```

```python
1  from time import time, sleep
2  from random import randint
3  from IPython.core.display import clear_output
4  from requests import get
5  from bs4 import BeautifulSoup
6  from warnings import warn
7
8  # Lists to store the scraped data
9  names = []
10 years = []
11 durations = []  # Adding duration to store movie durations
12 # IMDb ratings, Metascores, and votes are already present in the code
13
14 # Preparing the monitoring of the loop
15 start_time = time()
16 requests = 0
17
18 # Pages and years_url as per your previous code
19 pages = ['1', '2', '3', '4', '5']
20 years_url = ['2017', '2018', '2019', '2020']
21
22 for year_url in years_url:
23     for page in pages:
24         # Make a get request
25         response = get('https://hurawatch.cc/search/logan' + year_url +
26                        '&sort=num_votes,desc&page=' + page)
27
28         # Pause the loop
29         sleep(randint(8, 15))
30
31         # Monitor the requests
32         requests += 1
33         elapsed_time = time() - start_time
34         print('Request:{}; Frequency: {} requests/s'.format(requests, requests / elapsed_time))
35         clear_output(wait=True)
36
37         # Throw a warning for non-200 status codes
38         if response.status_code != 200:
39             warn('Request: {}; Status code: {}'.format(requests, response.status_code))
40
41         # Parse the content of the request with BeautifulSoup
42         page_html = BeautifulSoup(response.text, 'html.parser')
43
44         # Select all the movie containers from a single page
45         mv_containers = page_html.find_all('div', class_='fdi-item')
46
47         # For every movie in the containers
48         for container in mv_containers:
49             # Extracting movie name
50             name = container.h2.a.text
51             names.append(name)
52
53             # Extracting movie year
54             year = container.h2.find('span', class_='fdi-item').text
55             years.append(year)
56
57             # Extracting movie duration
58             duration = container.find('span', class_='fdi-item fdi-duration').text
59             durations.append(duration)
60
61 # Check the scraped data
62 print(names)
63 print(years)
64 print(durations)
65
```

```
-----------------------------------------------------------------------
KeyboardInterrupt                         Traceback (most recent call last)
<ipython-input-146-c6e20e626e83> in <cell line: 22>()
     23     for page in pages:
     24         # Make a get request
---> 25         response = get('https://hurawatch.cc/search/logan' + year_url +
     26                        '&sort=num_votes,desc&page=' + page)
     27
```

         ⇕ 13 frames

```
/usr/lib/python3.10/ssl.py in read(self, len, buffer)
   1157         try:
   1158             if buffer is not None:
-> 1159                 return self._sslobj.read(len, buffer)
   1160             else:
   1161                 return self._sslobj.read(len)

KeyboardInterrupt:
```
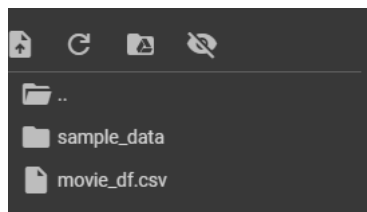
```
1 movie_ratings = pd.DataFrame({
2 'movie': names,
3 'year': years,
4 'duration': duration,
5 })
6 print(movie_ratings.info)
7 movie_ratings.head(10)
```

```
1 import pandas as pd
2 movie_df = pd.DataFrame({
3 'movie': names,
4 'year': years,
5 'duration': duration,
6 })
7 print(movie_df.info())
8 movie_df
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   movie     6 non-null      object
 1   year      6 non-null      object
 2   duration  6 non-null      object
dtypes: object(3)
memory usage: 272.0+ bytes
None
```

| | movie | year | duration |
|---|---|---|---|
| 0 | Logan | 2017 | 119m |
| 1 | Logan Lucky | 2017 | 119m |
| 2 | The Night Logan Woke Up | SS 1 | 119m |
| 3 | The Taking of Deborah Logan | 2014 | 119m |
| 4 | The Two Worlds of Jennie Logan | 1979 | 119m |
| 5 | Logan's Run | 1976 | 119m |

```
1 movie_df.to_csv('/content/movie_df.csv')
```



## ∨ Data preparation

```
1 movie_df['year'].unique()
```

```
array(['2017', 'SS 1', '2014', '1979', '1976'], dtype=object)
```